

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2001-005489

(43)Date of publication of application : 12. 01. 2001

(51)Int. Cl. G10L 15/22
G10L 15/00
G10L 15/28
H04N 5/00
H04Q 9/00

(21)Application number : 2000-117217 (71)Applicant : SONY INTERNATL EUROP
GMBH
(22)Date of filing : 13. 04. 2000 (72)Inventor : RAPP STEFAN
SILKE GORONJII
RALF KONPE
PETER BUCHNER
GIRON FRANCK
HELMUT LUCKE

(30)Priority

Priority	99 99107201	Priority	13. 04. 1999	Priority	EP
number :		date :		country :	

(54) CONTROL METHOD OF NETWORK

(57)Abstract:

PROBLEM TO BE SOLVED: To speedingly and flexibly control network devices by converting received and recognized voice commands into corresponding user network commands based on all pairs of user network commands included in user command interpretation elements and general languages.

SOLUTION: Vocabulary elements of first and second devices and corresponding commands L1 and L2 do not include same vocabulary elements. Therefore $L1 \cap L2 = \{\}$ and a received language L which is merged for interface description becomes $L = L1 \cup L2$. In other words a general document in a voice device 1 is constructed by combining pairs of commands corresponding to the vocabulary elements obtained from the device document which is described by the vocabulary elements of the

first device and the command L1 and the device document which is described by the vocabulary elements of the second device and the command L2. In other words user's uttered commands are converted into user network commands based on all pairs of user network commands included in user command interpretation elements and general languages.

CLAIMS

[Claim(s)]

[Claim 1] Change into a user network command characterized by comprising the following related in a user command in an audio station and by the above-mentioned related user network command. A control method of network equipment which controls network equipment connected to the above-mentioned network via a network.

A receiving step which receives at least one apparatus document including a language which corresponds to the above-mentioned network equipment and consists of at least one pair of a user network command relevant to the above-mentioned user command interpretation element.

An adaptation step which carries out adaptation of the apparatus document received [above-mentioned] to documents in general which consist of a language of this apparatus document and a language of the above-mentioned audio station constituted similarly.

A converting step which changes into a corresponding user network command a user's voice commanding which was above-received and has been recognized based on all the pairs of the above-mentioned user command interpretation element and a user network command included in the above-mentioned general purpose language.

[Claim 2] A step which the above-mentioned adaptation step judges for whether a language of the above-mentioned audio station and a language of a newly received apparatus document include at least one same user command interpretation element. When same user command interpretation element does not exist, update a language of the above-mentioned audio station and all the pairs of a user command interpretation element of a language of the above-mentioned audio station and a related user network command. A step combined so that all all the pairs of a user command interpretation element of a language of a newly [the above] received apparatus document and a related user network command may be contained. All the pairs of a user command interpretation element of a language of the above-mentioned audio station which updates a language

of the above-mentioned audio station and is not common when at least one same user command interpretation element exists and a related user network command join together and all the pairs of a user command interpretation element of a language of a newly [the above] received apparatus document and a related user network command. All the pairs of same user command interpretation element of a language of the above-mentioned audio station and a related user network command the above -- a control method of the network equipment according to claim 1 of having a step which gives an identifier which defines apparatus related to all the each of a pair of of a user command interpretation element of a language of a newly received apparatus document and a related user network command.

[Claim 3] A control method of the network equipment according to claim 2 wherein the above-mentioned identifier is a name of apparatus placed in front or postposed and given to a user command interpretation element of each set of a user network command relevant to the above-mentioned user command interpretation element.

[Claim 4] All the pairs of a user command interpretation element of a language of the above-mentioned audio station and a related user network command a step which combines all the pairs of a user command interpretation element of a language of a newly [the above] received apparatus document and a related user network command and which carries out and updates a language of the above-mentioned audio station [have and] each user command interpretation element -- a language of the above-mentioned audio station and the above -- a pair of user commands of all the of a user command interpretation element included in an intersection with a language of a newly received apparatus document and a related user network command. A control method of the network equipment according to claim 1 changing into a user network command corresponding based on all the pairs and selection processes of a user command interpretation element included in the above-mentioned general-purpose document and a related user network command.

[Claim 5] A control method of the network equipment according to claim 4 characterized by comprising the following.

A step which transmits an inquiry about which apparatus is used for the above-mentioned selection process to a user.

A step which receives and recognizes a reply from a user.

A step which chooses a corresponding user network command based on all the pairs of a related user network command included in a reply and a user command interpretation element and the above-mentioned documents in general from a user who has recognized.

[Claim 6]A control method of the network equipment according to claim 4 or 5 characterized by comprising the following.

A step to which the above-mentioned selection process transmits inquiry which apparatus was used immediately before to the above-mentioned network.

A step which receives and recognizes a reply from the above-mentioned network.

A step which chooses a corresponding user network command based on all the pairs of a user network command included in a reply and a user command interpretation elementand the above-mentioned documents in general from a network which received [above-mentioned].

[Claim 7]A control method of network equipment of claim 4 thru/or 6 given in any 1 paragraph characterized by comprising the following.

A step which the above-mentioned selection process receives the above-mentioned user's utteranceand is recognized.

A step which chooses a corresponding user network command based on all the pairs of utterance of a user who received and a user command interpretation elementand a user network command included in the above-mentioned documents in general.

[Claim 8]A control method of network equipment of claim 4 thru/or 7 given in any 1 paragraph characterized by comprising the following.

A step which transmits inquiry whether the above-mentioned selection process has a high possibility that which apparatus will be most used to the above-mentioned network.

A step which receives a reply from the above-mentioned network.

A step which chooses a corresponding user network command based on all the pairs of a related user network command included in a reply and a user command interpretation elementand the above-mentioned documents in general from the received above-mentioned network.

[Claim 9]A control method of the network equipment according to claim 8 that high network equipment of a possibility that ***** will also be used is chosen from network equipment in which a power supply is switched on.

[Claim 10]A control method of network equipment of claim 1 thru/or 9 characterized by transmitting an apparatus document corresponding to this network equipment to the above-mentioned audio station directly from predetermined apparatus after network equipment is connected to a

network given in any 1 paragraph.

[Claim 11]A control method of network equipment of claim 1 thru/or 10 characterized by transmitting an apparatus document to the above-mentioned audio station from apparatus predetermined [this] after the above-mentioned audio station transmits a user network command to predetermined apparatus given in any 1 paragraph.

[Claim 12]. [whether network equipment corresponding to the above-mentioned apparatus document is connected to the above-mentioned networkand] Or a control method of network equipment of claim 1 thru/or 11 which will be characterized by transmitting this apparatus document to the above-mentioned audio station from this apparatus document providing device if a user network command which requires an apparatus document with an apparatus document providing device specific from the above-mentioned audio station is received given in any 1 paragraph.

[Claim 13]A control method of network equipment of claim 1 thru/or 12wherein a general-purpose document stored in the above-mentioned audio station is empty in first stage given in any 1 paragraph.

[Claim 14]A control method of network equipment of claim 1 thru/or 12wherein a general-purpose document stored in the above-mentioned audio station includes a pair of basic set of a user command interpretation element and a related user network command in first stage given in any 1 paragraph.

[Claim 15]A control method of the network equipment according to claim 14wherein the above-mentioned user command interpretation element and a pair of corresponding basic set of a user network command define initial grammar about an expression of utterance.

[Claim 16]A control method of network equipment of claim 1 thru/or 15wherein the above-mentioned user command interpretation element contains a lexical element given in any 1 paragraph.

[Claim 17]A control method of network equipment of claim 1 thru/or 16wherein the above-mentioned user command interpretation element includes a definition of grammar about continuous utterance of a user who may appear based on a keyword and/or a category given in any 1 paragraph.

[Claim 18]A control method of network equipment of claim 1 thru/or 17wherein the above-mentioned user command interpretation element includes a definition about pronunciation given in any 1 paragraph.

[Claim 19]A control method of the network equipment according to claim 18wherein a definition about the above-mentioned pronunciation is related with a language of arbitrary apparatus.

[Claim 20]A control method of the network equipment according to claim

18wherein a definition about the above-mentioned pronunciation is related with a language of two or more apparatus.

[Claim 21]A control method of network equipment of claim 1 thru/or 20wherein the above-mentioned user command interpretation element includes at least one word sequence given in any 1 paragraph.

[Claim 22]The above-mentioned user command interpretation element is the control method of network equipment given in 1 paragraph to either of claims 1 thru/or 21 including a rule which defines that of arbitrary user commandsand mapping between concepts.

[Claim 23]A control method of network equipment of claim 1 thru/or 22wherein the above-mentioned user command interpretation element includes a user command interpretation element for other productive wordsand the same information as a related pair of a user network command given in any 1 paragraph.

[Claim 24]A control method of network equipment given in any 1 paragraph of claims 1 thru/or 23wherein the above-mentioned user command interpretation element includes standardized information about a function of a category of apparatusand/or this apparatus.

[Claim 25]A language of network equipmentor a language of the above-mentioned audio stationThis above-mentioned audio station to one received user command including a reasoning part by this reasoning part. A control method of network equipment given in the any 1 paragraph according to claim 1 to 24 transmitting several different user network commands to 2 or two or more network equipment based on information about a category and/or a function of network equipment.

DETAILED DESCRIPTION

[Detailed Description of the Invention]

[0001]

[Field of the Invention]This invention relates to the voice interface which can extend a vocabulary dynamically and actively in home network environment based on the vocabulary depending on the apparatus or the medium transmitted from network equipment. Especially this invention relates to the control method of the network equipment which can extend a vocabulary in the audio station which realizes a voice interface etc. Hardwaressuch as a videotape recorder (VTR)may constitute the device which applied this inventionfor exampleor softwaresuch as an electronic program guide (electronic programming guide)may be sufficient as itfor example.

[0002]

[Description of the Prior Art] In home network environment the apparatus which transmits the vocabulary and voice interface which describe their function to an audio station is indicated by European patent publication-before-examination EP-A-97118470 No. An audio station changes into a corresponding user network command a user's utterance received and recognized and controls this apparatus based on this user network command.

[0003] Motorola (Motorola) exhibited the language relevant to the Vox markup language 1.0 (VoxML1.0) based on extensible markup language XML (below Extensible Markup Language: calls it XML.) in September 1998. This language is used for describing a dialog (dialog) system by specifying a prompt and the dialog step which generally consists of a list of available options. By hypermedia description language HTML (below Hypertext Markup Language: calls it HTML.) with a rich text. According to VoxML1.0 description of voice application becomes easy so that a picture a hypertext link and graphic user interface input control can be described easily. This VoxML1.0 includes the information about the pointer which directs other references which are useful when developing the composition of the syntax of the attribute of two or more elements and these elements an example a Vox markup language document or a dialog and the application which uses a Vox markup language.

[0004] Similarly by the Euro speech 97 Europe Acoustical Society of America-sponsored [in the low toss island in Greece]. "The paper announced by Richard SUPURTO and others The markup language for text-to-speech synthesis ISSN. 1774 pages of 1018-4074. "a markup language for text-to-speech synthesis" by Richard Sproat et al. ESCA. Eurospeech 97 Rhodes [Greece and] The voice text markup language (spoken text markup language) (STML) is indicated by ISSN 1018-4074 and page 1774. This STML provides a text voice (TTS) composing device (text-to-speech synthesizers) with the knowledge about the composition of a text. For example in a multilingual TTS system a STML text sets up language initializes a speaker to the language and a suitable language and a speaker specific table are loaded.

[0005] Philips (Phillips) developed the dialog description language called HDDL. HDDL is a language for the dialog especially used in an automatic inquiry system. A dialogue system is sold after being created by offline mode using HDDL.

[0006]

[Problem(s) to be Solved by the Invention] Drawing 3 is a figure showing the conventional audio station 21. The audio station 21 is connected to

the microphone 22, the loudspeaker 30 and the bus 40. After the input signal of the microphone 22 is processed by the digital signal processor (DSP) 23 which is built in the memory 23, it is supplied to the central processing unit (CPU) 24. The result of an operation is outputted to the loudspeaker 30 via DSP 23 which is built in the memory 23, or CPU 24 outputs it to the bus 40 via the link layer control section 25 and the I/F physical layer part 26. DSP 23 and CPU 24 can also access the memory 28. All the information required in order to control processing of an input signal is stored in the memory 28. DSP 23 accesses the feature extraction part 27e of the memory 27. All the information required for speech recognition and voice synthesis is stored in the memory 27. CPU 24 also accesses not only each voice interface definition part and general-purpose voice interface definition part for apparatus #1, apparatus #2 and apparatus #3 but a recognition part and the secretary matter / phoneme converter 27f, two or more apparatus and here. Each voice interface definition part and a general-purpose voice interface definition part are independently memorized in the memory 27.

[0007] Thus, since it had the general-purpose voice interface definition part separately while providing the voice interface definition part according to apparatus, the command which two or more apparatus shares needed to be managed individually, structure became complicated and the efficiency of processing or analysis was bad in the conventional audio station 21.

[0008] Then, this invention is made in view of the actual condition mentioned above and is a thing.

The purpose is to provide the control method of simple and efficient network equipment of transmitting the function of ** and network equipment and a voice interface to a network audio station and operating the function of two or more network equipment and a voice interface in an audio station.

[0009] That is, the purpose of this invention is [for controlling within a network the network equipment connected to the network by an audio station] easy, is quick and is providing the method of being supplied. An audio station is a device which changes a user command into a user network command and controls network equipment via a network based on the function of network equipment and a voice interface.

[0010]

[Means for Solving the Problem] According to this invention, all the network equipment connected to a network is equivalent to at least one apparatus document which defines a function of network equipment and a

voice interface. Network equipment itself may be provided with an apparatus document or it may be drawn up to the exterior of network equipment. a pair with a user network command relevant to a user command interpretation element of network equipment in an apparatus document -- 1 set -- or it has two or more sets. An audio station receives an apparatus document, unifies a received apparatus document and creates one voice interface description. This voice interface description integrated is created based on a language of an audio station. Voice interface description can be referred to as a general-purpose document for all the networks. A user's utterance command received and recognized by audio station is changed into a user network command based on all the pairs of a user command interpretation element and a related user network command included in a general purpose language. A user command interpretation and an execution element have a definition of a lexical element and grammar and a pronunciation definition for example.

[0011] the [European patent gazette] -- an audio station is fetched from a database which exists in memory storage or a remote which an audio station equips with an apparatus document as indicated by EP-A-97118470 No.

[0012] After an apparatus document is received from network equipment at the time of execution, adaptation of the general-purpose document may be purely carried out syntactically in an audio station. Network equipment may have two or more apparatus documents. Each of those documents describe a part of function of network equipment and some interfaces. And among those only a document actually needed is transmitted to an audio station. Such a document is a document which defines a part of function of a document or network equipment which described a function of network equipment using a certain language for example. When it judges with an audio station or each network equipment itself needing further phonetic function of network equipment of these each an apparatus document corresponding to the further phonetic function can be transmitted to an audio station. An audio station carries out adaptation of the general-purpose document based on this further document and generates a user network command corresponding based on that general-purpose document by which adaptation was carried out at the time of execution.

[0013]

[Embodiment of the Invention] Drawing 1 is a figure showing the audio station 1 which applied this invention. The audio station 1 is connected to the microphone 2 and the loudspeaker 10 and the bus 20. After the input signal based on the sound inputted into the microphone 2 is processed by the digital signal processor (DSP) 3 which built in the

memory 3ait is supplied to the central processing unit (CPU) 4. The result of an operation is outputted to the loudspeaker 10 via DSP9 which built in the memory 9aor CPU4 outputs it to the bus 20 via the link layer control section 5 and the I/F physical layer part 6. DSP3 and CPU4 can also access the memory 8. All the information required in order to control processing of an input signal is stored in the memory 8. DSP3 accesses the feature extraction part 7e of the memory 7. All the information required for speech recognition and voice synthesis is stored in the memory 7.

[0014]In the conventional audio stationtwo or more voice interface definition parts corresponding to each of two or more connected network equipment and one general-purpose voice interface definition part were provided. On the other handthe audio station 1 which applied this invention has the only voice interface definition part integrated. This voice interface definition part integrated supports the general-purpose document described by the conventional general-purpose voice interface definition part.

[0015]Drawing 2 is a functional-blocks figure of the network equipment 11 which applied this invention. The network equipment 11 is provided with CPU12. CPU12 carries out an interaction to the software 15 for appliance control for controlling the network equipment 11the memory 14the link layer control section 17and the I/F physical layer part 16and outputs various information to the bus 10. The interaction of CPU12 can also be carried out to the memory 13. According to this inventionthe memory 13 has memorized a voice interface definitioni.e. lor two or more apparatus documents.

[0016]As mentioned abovethe network equipment 11 does not necessarily need to be provided with the memory 13 which has a voice interface definition part. The voice interface definition of the network equipment 11 may be provided with a voice interface definition providing device. The audio station 1 shown in drawing 1 can be accessed at a voice interface definition providing device.

[0017]The two network equipment 11 is connected to the network provided with the one audio station 1 in the example of this invention explained below. Two apparatus documents in which the general-purpose document of the audio station 1 defines the language of the two network equipment 11 by empty are merged into description of one interface in the audio station 1 in first stage (merge). The number of the network equipment 11 connected to a network may be [two or more / arbitrary] n. Based on the newly received apparatus documentadaptation of the general-purpose document of the audio station 1 can be carried out. Belowin order to

explain briefly an apparatus document shall have a user command interpretation element which consists of one lexical element.

[0018] In the following explanation L1 shows the lexical element and the corresponding command of an acceptance language (accepted language) i.e. the 1st apparatus. L2 shows the 2nd lexical element and corresponding command of apparatus similarly. According to mathematical expression L1 is a set of at least one lexical element i.e. the user network command corresponding to word w_i and this word w_i . Word w_i may not be a single word but may be perfect utterance which consists of two or more words. An acceptance language may be added to a lexical element and may also include the rule for the grammar for pronunciation and a word sequence and/or speech comprehension (speech understanding) and a dialog for example.

[0019] In the 1st example L1 and L2 do not contain the same lexical element i.e. a common word. Therefore the acceptance language L which is $L1 * L2 = \{\}$ and was merged for interface description is set to $L = L1 * L2$. That is the general-purpose document in the audio station 1 is built by annexing the pair of the lexical element obtained from the apparatus document 1 described by the language L1 of the 1st apparatus and the apparatus document 2 described by the language L2 of the 2nd apparatus and a corresponding command. Since L1 and L2 do not contain the same lexical element in order to suggest a thing [as opposed to / command / corresponding / which apparatus in the lexical element] a lexical element is a user network command is generated appropriately and is correctly transmitted to corresponding apparatus.

[0020] In the 1st example two network equipment is a television set and a CD player for example. At this time L1 corresponding to a television set and L2 corresponding to a CD player consist of the following lexical elements within an apparatus document respectively.

[0021] $L1 = \{\text{MTVCNN}\}$

$L2 = \{\text{reproduction and a stop}\}$

Since L1 and L2 do not contain the same lexical element i.e. it is $L1 * L2 = \{\}$ the acceptance language L with which interface description was merged becomes $L = L1 * L2 = \{\text{MTVCNN reproduction and a stop}\}$. For example these lexical elements correspond with the user network command which has a function which the following lists are shown respectively.

[0022]

> reproduction which switches CNN [<switches television to MTV>] -> <television to CNN -> carrying out only > stop -> <a CD player is made into stop mode> which makes < CD player reproduction mode. [MTV->]
When sharing the lexical element with two same apparatus i.e. the same

word it becomes $L1**L2 \neq \{\}$. According to this invention it is identifiable also in a vocabulary common to these. The name of apparatus is prefaced or postposed in the 1st example of this invention by the same lexical element between each acceptance language (i.e. between $L1$ and $L2$) that forms an intersection at least. therefore a user -- before the name of each apparatus of a request of a user's utterance command -- and/or it must add back. As mentioned above when a command is not common it is not necessary to add the name of apparatus but and it may add.

[0023] In the following example [2nd] the name of apparatus is added before each command and the new language L which forms description of an interface is shown in the language of the apparatus by which annexation of a word without a possibility that it may be mixed up between the two languages $L1$ and $L2$ and the name of apparatus were prefaced. It will be set to $L=L1**(L1**L2) **L2**(L1**L2) **n1L1**n2L2$ if the name of apparatus is set to $n1$ and $n2$ respectively.

[0024] Below an above-mentioned example is described to an example for a network provided with a CD reproduction device and a tape reproducer and the control method of the network equipment by this invention is clarified. Here the CD reproduction device and the tape reproducer are named "CD" and a "tape" respectively. The lexical element which constitutes the acceptance language $L1$ of a CD reproduction device and the acceptance language $L2$ of a tape reproducer respectively is shown below.

[0025] $L1 = \{\text{reproduction a stop and a skip}\}$

$L2 = \{\text{playback sound recording a stop and rewinding}\}$

The acceptance language of a voice interface is $L = (L1**(L1**L2) ** (it becomes L2**(L1**L2) **n1L1**n2L2 = \{\text{a skip sound recording rewinding CD reproduction CD stop CD skip tape recording a tape stop and tape rewinding}\})$. The function of a user network command to correspond is as follows.

[0026]

Skip -> > sound-recording -> [</ which is skipped to the next track of CD] <a tape reproducer is made into sound recording mode> Rewind and it skips to the next track of CD skip [<suspends playback of a CD reproduction device>] -> <CD. [CD stop / <makes a CD reproduction device reproduction mode> / ->] [CD reproduction / <rewinds a tape> / ->] [->] The tape stop to carry out and <makes a tape reproducer sound recording mode> -> tape rewinding [<which suspends playback/sound recording of a tape reproducer>] -> <a tape is rewound> [tape recording -> / <makes a tape reproducer reproduction mode>] [> tape reproduction ->] [0027] In this 2nd example when the recognized command

is equivocal the same word problem (same words problem) asks automatically to a user and is solved by making the apparatus for which it asks to a user specify. Although the scenario in this case is the same as the language which has the distinctiveness in the 1st example formally that interpretation changes. The 3rd example that clarifies this example that specifies the variation of this interpretation is based on the scenario which is based on the 2nd example namely was given to the CD reproduction device and tape reproducer as network equipment. The acceptance language L is formed like the acceptance language L of the 1st example in that a selection process is performed when recognizing the lexical element contained in the common area of the acceptance language L1 and the acceptance language L2. That is it is $L=L1**L2$. The lexical element which constitutes the acceptance language L of an audio station under the conditions of having the lexical element with respectively same the acceptance language L1 and the acceptance language L2 like the 2nd example becomes $L = \{\text{playback a stop a skip sound recording and rewinding}\}$. the user network command to which it corresponds in this case is a playback \rightarrow < clear imitation -- it is a #1> stop \rightarrow < clear imitation -- a sound recording [<is skipped to the next track of CD>] \rightarrow < tape reproducer is made into sound recording mode -- it > rewinds and is \rightarrow <a tape is rewound>. [#2> skip \rightarrow] here -- " -- clear -- an imitation -- #1" means that a voice recognition part outputs an inquiry called the command or ("on what device") to which apparatus and switches the acceptance language L to acceptance language La. La consists of the lexical element $La = \{\text{CD and tape}\}$ and the function of a user network command to correspond It is CD \rightarrow <a CD reproduction device's is made into reproduction mode> and <switchbacking to L> and is tape \rightarrow <a tape reproducer's is made into reproduction mode> and <switchbacking to L>.

[0028] here -- " -- clear -- an imitation -- #2" means that the audio station 1 outputs the inquiry to which apparatus whether it is a command and switches the acceptance language L to the acceptance language Lb. Lb consists of a lexical element indicated to be $Lb=La = \{\text{CD and tape}\}$ and the function of a user network command to correspond It is CD \rightarrow <a CD reproduction device's is made into stop mode> and <switchbacking to L> and is Tape \rightarrow <a tape reproducer's is made into stop mode> and <switchbacking to L>.

[0029] In the 3rd example of this invention the control management of the audio station shown in the 2nd example can also be shortened. In this case the audio station 1 considers that the apparatus used immediately before is default apparatus and changes into a corresponding user network

command a user's utterance command received and recognized based on the general-purpose document which consists of the acceptance language L with which the audio station 1 and default apparatus were unified. Also when controlling again other apparatus other than the apparatus used immediately before the control management of an audio station can be shortened. In this case it directs to control apparatus to an audio station by only emitting the name of apparatus for example.

[0030] Adaptation of a dialog is performed in the 4th example. In the above situations in order to assign an equivocal word correctly an inquiry of apparatus clear statement is outputted to a user and a system recognizes a user's reaction to the inquiry. When the great portion of reply of the user to the inquiry for apparatus clear statement is the apparatus specified immediately before an inquiry is skipped and it may be made to carry out direct supply of the user network command to last apparatus. In order to predict to which apparatus the command was taken out two or more sources of information may be used. For example it not only takes into consideration the apparatus used immediately before but it takes into consideration the points of comparison of each apparatus or a classification of apparatus. For example a record command has a high possibility of being not the record to an audiotape but record on videotape when being taken out while the user is watching television. The record command taken out on the other hand while the user is listening to radio has a high possibility of being record to an audiotape. Then it may be made to choose the high network equipment of a possibility of being used from the network equipment in which the power supply is switched on for example. A classification of apparatus can also be learned by investigating the description about a user's act or the function of apparatus.

[0031] In all the above-mentioned examples after as for the rear stirrup by which network equipment was connected to the network apparatus receives a user network command from the audio station 1 each network equipment can transmit an apparatus document to the audio station 1 directly. The apparatus document etc. which became independent of the present status may be transmitted to an audio station. Apparatus may have two or more apparatus documents. Or apparatus may update a dialog and voice recognition capability when changing a function dynamically based on change of contents information (content information) and accepting the change. Or a function may be periodically changed by transmitting a new document to an audio station. An apparatus document may be stored for example in an apparatus document providing device like a device manufacturer's internet server without being stored in corresponding

apparatus. In this case if the audio station 1 recognizes it as the above-mentioned apparatus being connected to a network it may download an apparatus document.

[0032] Apparatus can have two or more documents including the acceptance language of the whole apparatus by the language of the country where each differs. For example such apparatus transmits a German apparatus document to an audio station first. Then after receiving a command corresponding from a user an English apparatus document can be transmitted to an audio station. Thereby the audio station can change the user command of either German or English into a user network command and can control apparatus.

[0033] The definition of the pronunciation for the word used when performing recognition or composition other than a lexical element may be included as an element of an apparatus document. And the definition of those pronunciation may be directly unified in the general-purpose document of an audio station. When the definition of pronunciation is included in language the definition of those pronunciation may be used for one recognition for apparatus and a synchronizer for exclusive use or all the apparatus may share it for example. That is the definition of the pronunciation introduced by the 1st apparatus can be validated also to the word emitted to the 2nd apparatus i.e. other apparatus connected to the network.

[0034] In order to compound two or more apparatus documents may be unified and adaptation of a general-purpose document may be performed in an audio station. The reply to an inquiry may be performed not only to one apparatus but to two or more apparatus. It is necessary to guarantee that composition of the utterance from two or more apparatus is not outputted simultaneously. Some utterance is produced according to the phenomenon of the exterior like warning based on a supply priority. For example after those utterance interrupts other utterance and a user's inactive period expires the interrupted utterance is reintroduced and a dialog continues. While the stream considered that the likelihood which is in agreement with a user's input and a user's utterance is high is continuing the stream of two or more dialogs is parallel and is managed on the assumption that the fact that all the streams are in an active state.

[0035] The general-purpose document of an audio station may have a basic set of the element which may be empty in first stage or describes the interface for 1 or two or more network equipment. In the case of the latter the document transmitted to an audio station from apparatus does not have concreteness but refer to the general-purpose document for it selectively.

[0036]For exampleonly the dialog grammar for specific apparatus may be specified in an apparatus documentand the default grammar for the expression of utterance may be included in the fundamental set of the user command interpretation element of a general-purpose document. A certain function may be thoroughly described in the general-purpose document of a voice recognition part. For exampleit is described how an electronic program guide is controlled by a voice dialog. In this casethe variable which shows the information about a program namea maker namean actor namea time zoneand a broadcasting day is included in default grammar and the grammar for a series of wordsfor exampleand the document transmitted to an audio station from an electronic program guide device includes only the information by which these variables are fulfilled.

[0037]In other examplesthe apparatus document transmitted to an audio station from apparatus contains the keyword and category identifier corresponding to a network command. An audio station determines which grammar about a word sequence is applied to utterance of the continuing user who may be generated based on these category identifiers. For examplewhen a keyword is "Voice of America (voice of America)" and a category is "radio"a userIt generates saying "I liking to hear Voice of America (I want to listen to voice of America)" or "one [Voice of America] (Turn on voiceof America)." When a category is a "radio cassette recorder recorder"utterance of the user who may be generated is "recording Voice of America (Please record voice of America)" etc.

[0038]required (time8:00) in order to start the group of the concept/value which forms the foundation of speech comprehension in order to perform adaptation of a general-purpose documentfor examplea certain processing-- it may be made to contain in the apparatus document to which the said group is transmitted by the audio station as grammar for apparatus As for this rulethe apparatus document which describes a voice interface defines [mapping between a word sequence and a conceptand] how a variable value is buried including a rule. An apparatus document has a rule which defines mapping to the operation from the group of a couple for a concept/value. For exampleit goes away with a (command and record) (title and a windand a group called **)and (time and 8:00) expresses the picture recording processing of VTR. A dialog may be begun from the user sidenamelythe turn of utterance is good also considering a user as initiative. "8 o'clockthe turn of utterance is left with a movie and a wind and is **" etc.for example. When information is given by 1 or two or more utterancethe turn of utterance is left to a user. A system collects information required in

order to put operation into operation for example by the user network command transmitted to each network equipment based on grammar. A user's utterance which this system transmitted the further document that includes additional grammar for example to the audio station the audio station could ask back information to the user when information was missing and has been recognized newly can also be assigned to the further user network command.

[0039] An apparatus document as a user command interpretation element Grammar for a lexical element pronunciation and a word sequence Besides the information about the rule for speech comprehension and a dialog the information for the grammar for the same information corresponding to two or more languages assigned to a same or equivalent user network command i.e. a lexical element pronunciation and a word sequence speech comprehension and a dialog which carries out rule Seki may also be included. The user can control apparatus by this information with any languages. As mentioned above this information may be included in two or more apparatus documents which the same apparatus has. The interface of one audio station which can realize processing with the sound of two or more languages by this can be specified. Since the pair of the above-mentioned concept/value is language independence grammar of a dialog can be considered as language independence.

[0040] As mentioned above an apparatus document contains the vocabulary given clearly or suggestively by the grammar for utterance. 1 or two or more pronunciation can also be arbitrarily given to each lexical element i.e. each word. When arbitrary pronunciation is not given to an apparatus document this may be generated automatically. However when the special pronunciation in the case where especially words are a proper noun a foreign word and an abbreviation a dialect or a foreign language must be recognized in this case it is easy to generate an error. Pronunciation serves as the foundation for generating the model of a word in the recognition part to which the set of the model of a single sound or an allophone was given.

[0041] An apparatus document besides the user command interpretation element which creates the compound command about two or more available apparatus and a related user command The information about functions like "the standardized information for example the information about the category of apparatus like "VTR" and it has a recordable tape" may also be included. For example a user emits saying "record this movie (Please record this movie)" when choosing the movie of the schedule broadcast in the future from the electronic program guide displayed on a television set. That is a network command needs to transmit the information about a

suitable channels and a time zone to VTR from a program guide device. And it is necessary to transmit the command which guarantees that suitable recording is performed to VTR. As a situation similar to this a movie may be dubbed from a certain VTR tape to another VTR tape. The voice interface description of each apparatus does not usually define the utterance about such utterance, i.e. a suitable channels and a time zone etc. by such a situation. Therefore an audio station may carry out adaptation of the general-purpose document based on reasoning by a reasoning part. The reasoning part is contained in the apparatus document as a reasoning component and provides the additional function which may be included in one of the compound commands about two or more apparatus with description of the voice interface containing the grammar and the dialog about a word sequence.

[0042] Since all the information for controlling network equipment is included in one general-purpose document in an audio station it can conduct processing or analysis easily. It is easy especially when [especially] relation is between the double information assigned to the same or equivalent user network command between the available contents from which a user command interpretation element differs or/and in a general-purpose document.

[0043] Therefore according to this invention network equipment controllable via the audio station with which a network is provided can transmit the apparatus document which describes the function of apparatus and a voice interface to an audio station. An audio station unifies an apparatus document in a general-purpose document. A general-purpose document forms the foundation for changing the recognized user command into a user network command using a user command interpretation element in order to control the connected network equipment. An apparatus document consists of a vocabulary like for example a user command interpretation element and a related user network command. The rule for the grammar not only for the same information as these corresponding to two or more languages or the information about the dynamic dialog in speech comprehension but pronunciation and a word sequence speech comprehension and a dialog may be included in the user command interpretation element of an apparatus document. One apparatus has two or more apparatus documents and may transmit them to an audio station dynamically at the time of necessity. Since network equipment transmits the specification about a phonetic function to an audio station at the time of execution apparatus can change the function dynamically by applying this invention based on change of the contents of the apparatus document. When apparatus changes the function dynamically network equipment generates an apparatus

document dynamically or updates the existing apparatus document. Updating is performed by updating or inserting the name of a broadcasting station for example and network equipment transmits the updated apparatus document to an audio station.

[0044]

[Effect of the Invention] As mentioned above in the control method of the network equipment concerning this invention. Correspond to network equipment and at least one apparatus document including the language which consists of at least one pair of the user network command relevant to a user command interpretation element is received adaptation of the received apparatus document is carried out to the general-purpose document which consists of a language of an apparatus document and a language of the audio station constituted similarly -- adaptation being carried out and A user's voice commanding received and recognized is changed into a corresponding user network command based on all the pairs of a user command interpretation element and the related user network command included in a general purpose language. Thereby according to the control method of the network equipment concerning this invention the user network command which controls two or more network equipment connected to the network provided with an audio station can be managed processed and analyzed simply and efficiently.

DESCRIPTION OF DRAWINGS

[Brief Description of the Drawings]

[Drawing 1] It is a block diagram showing the composition of the audio station which applied this invention.

[Drawing 2] It is a block diagram showing the composition of the network equipment which applied this invention.

[Drawing 3] It is a block diagram showing the composition of the conventional audio station.

[Description of Notations]

1 An audio station and 2 [Bus] A microphone 3 DSP 4 CPU 5 link-layer control section 6 I/F physical layer part and 7 A memory and 8 A memory 9 DSP ten loudspeakers and 20

【特許請求の範囲】

【請求項1】 音声装置においてユーザコマンドに関連するユーザネットワークコマンドに変換し、上記関連するユーザネットワークコマンドにより、ネットワークを介して上記ネットワークに接続されたネットワーク機器を制御するネットワーク機器の制御方法において、上記ネットワーク機器に対応し、上記ユーザコマンド解釈要素と関連するユーザネットワークコマンドの少なくとも1つの対からなる言語を含む少なくとも1つの機器文書を受信する受信ステップと、上記受信された機器文書を、該機器文書の言語と同様に構成される上記音声装置の言語からなる汎用文書に適応化する適応化ステップと、上記受信及び認識されたユーザの音声コマンドを、上記ユーザコマンド解釈要素と上記汎用言語に含まれるユーザネットワークコマンドの全ての対に基づいて、対応するユーザネットワークコマンドに変換する変換ステップとを有するネットワーク機器の制御方法。

【請求項2】 上記適応化ステップは、上記音声装置の言語と新たに受信された機器文書の言語が少なくとも1つの同様なユーザコマンド解釈要素を含んでいるか否かを判定するステップと、同様なユーザコマンド解釈要素が存在しない場合、上記音声装置の言語を更新し、上記音声装置の言語のユーザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対と、上記の新たに受信された機器文書の言語のユーザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対とが全て含まれるように結合するステップと、少なくとも1つの同様なユーザコマンド解釈要素が存在する場合、上記音声装置の言語を更新し、共通しない上記音声装置の言語のユーザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対と、上記の新たに受信された機器文書の言語のユーザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対とを結合し、同様な上記音声装置の言語のユーザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対と、上記新たに受信された機器文書の言語のユーザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対のそれぞれに対して関連する機器を定義する識別子を付与するステップとを有する請求項1記載のネットワーク機器の制御方法。

【請求項3】 上記識別子は、上記ユーザコマンド解釈要素と関連するユーザネットワークコマンドの各対のユーザコマンド解釈要素に対して前置又は後置されて付与される機器の名称であることを特徴とする請求項2記載のネットワーク機器の制御方法。

【請求項4】 上記音声装置の言語のユーザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対と、上記の新たに受信された機器文書の言語のユー

ザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対とを結合するして上記音声装置の言語を更新するステップを有し、

各ユーザコマンド解釈要素が上記音声装置の言語と上記新たに受信された機器文書の言語との共通部分に含まれているユーザコマンド解釈要素及び関連するユーザネットワークコマンドの対の全てのユーザコマンドは、上記汎用文書に含まれるユーザコマンド解釈要素及び関連するユーザネットワークコマンドの全ての対と選択処理とに基づいて対応するユーザネットワークコマンドに変換されることを特徴とする請求項1記載のネットワーク機器の制御方法。

【請求項5】 上記選択処理は、ユーザに対してどの機器を使用するかに関する問い合わせを送信するステップと、ユーザからの回答を受信し認識するステップと、認識したユーザからの回答及びユーザコマンド解釈要素と上記汎用文書に含まれる関連するユーザネットワークコマンドの全ての対に基づいて、対応するユーザネットワークコマンドを選択するステップとを有する請求項4記載のネットワーク機器の制御方法。

【請求項6】 上記選択処理は、上記ネットワークに対してどの機器が直前に使用されたかという問い合わせを送信するステップと、上記ネットワークからの回答を受信し認識するステップと、上記受信したネットワークからの回答及びユーザコマンド解釈要素と上記汎用文書に含まれるユーザネットワークコマンドの全ての対に基づいて、対応するユーザネットワークコマンドを選択するステップとを有する請求項4又は5記載のネットワーク機器の制御方法。

【請求項7】 上記選択処理は、上記ユーザの発話を受信し認識するステップと、受信したユーザの発話及びユーザコマンド解釈要素と上記汎用文書に含まれるユーザネットワークコマンドの全ての対に基づいて、対応するユーザネットワークコマンドを選択するステップとを有する請求項4乃至6のいずれか1項記載のネットワーク機器の制御方法。

【請求項8】 上記選択処理は、上記ネットワークに対してどの機器が最も用いられる可能性が高いかという問い合わせを送信するステップと、上記ネットワークからの回答を受信するステップと、受信した上記ネットワークからの回答及びユーザコマンド解釈要素と上記汎用文書に含まれる関連するユーザネットワークコマンドの全ての対に基づいて、対応するユーザネットワークコマンドを選択するステップとを有する請求項4乃至7のいずれか1項記載のネットワーク機器の制御方法。

【請求項9】 上記最も用いられる可能性の高いネットワーク機器は、電源が投入されているネットワーク機器

から選択される請求項8記載のネットワーク機器の制御方法。

【請求項10】 ネットワーク機器がネットワークに接続された後に、該ネットワーク機器に対応する機器文書が所定の機器から上記音声装置に直接送信されることを特徴とする請求項1乃至9のいずれか1項記載のネットワーク機器の制御方法。

【請求項11】 上記音声装置が所定の機器にユーザネットワークコマンドを送信した後に、該所定の機器から機器文書が上記音声装置に送信されることを特徴とする請求項1乃至10のいずれか1項記載のネットワーク機器の制御方法。

【請求項12】 上記機器文書に対応するネットワーク機器が上記ネットワークに接続されるか、あるいは、機器文書提供装置が上記音声装置から特定の機器文書を要求するユーザネットワークコマンドを受信すると、該機器文書提供装置から上記音声装置に該機器文書が送信されることを特徴とする請求項1乃至11のいずれか1項記載のネットワーク機器の制御方法。

【請求項13】 上記音声装置内に格納されている汎用文書は、初期的には空であることを特徴とする請求項1乃至12のいずれか1項記載のネットワーク機器の制御方法。

【請求項14】 上記音声装置内に格納されている汎用文書は、ユーザコマンド解釈要素及び関連するユーザネットワークコマンドの対の基本セットを初期的に含んでいることを特徴とする請求項1乃至12のいずれか1項記載のネットワーク機器の制御方法。

【請求項15】 上記ユーザコマンド解釈要素及び対応するユーザネットワークコマンドの対の基本セットは、発話の言い回しに関する初期的文法を定義することを特徴とする請求項14記載のネットワーク機器の制御方法。

【請求項16】 上記ユーザコマンド解釈要素は、語彙要素を含むことを特徴とする請求項1乃至15のいずれか1項記載のネットワーク機器の制御方法。

【請求項17】 上記ユーザコマンド解釈要素は、キーワード及び／又はカテゴリに基づいて、出現する可能性のあるユーザの連続的発話に関する文法の定義を含むことを特徴とする請求項1乃至16のいずれか1項記載のネットワーク機器の制御方法。

【請求項18】 上記ユーザコマンド解釈要素は、発音に関する定義を含むことを特徴とする請求項1乃至17のいずれか1項記載のネットワーク機器の制御方法。

【請求項19】 上記発音に関する定義は、任意の機器の言語に関連付けられていることを特徴とする請求項18記載のネットワーク機器の制御方法。

【請求項20】 上記発音に関する定義は、複数の機器の言語に関連付けられていることを特徴とする請求項19記載のネットワーク機器の制御方法。

【請求項21】 上記ユーザコマンド解釈要素は、少なくとも1つの単語列を含むことを特徴とする請求項1乃至20のいずれか1項記載のネットワーク機器の制御方法。

【請求項22】 上記ユーザコマンド解釈要素は、任意のユーザコマンドの概念の間のマッピングを定義する規則を含むことを特徴とする請求項1乃至21のいずれかに1項記載のネットワーク機器の制御方法。

【請求項23】 上記ユーザコマンド解釈要素は、他の発話言語用のユーザコマンド解釈要素及び関連するユーザネットワークコマンドの対と同様の情報を含むことを特徴とする請求項1乃至22のいずれか1項記載のネットワーク機器の制御方法。

【請求項24】 上記ユーザコマンド解釈要素は、機器のカテゴリ及び／又は該機器の機能に関する標準化された情報を含むことを特徴とする請求項1乃至23のいずれかに1項に記載のネットワーク機器の制御方法。

【請求項25】 ネットワーク機器の言語又は上記音声装置の言語は、推論部を含み、該上記音声装置は、1つの受信したユーザコマンドに対し、該推論部により、ネットワーク機器のカテゴリ及び／又は機能に関する情報に基づいて、2又は複数のネットワーク機器に複数の異なるユーザネットワークコマンドを送信することを特徴とする請求項1乃至24記載のいずれか1項に記載のネットワーク機器の制御方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、ホームネットワーク環境において、ネットワーク機器から送信される機器又は媒体に依存した語彙に基づいて、語彙を動的にそして能動的に拡張することのできる音声インターフェイスに関する。特に、本発明は、音声インターフェイス等を実現する音声装置内で、語彙を拡張できるネットワーク機器の制御方法に関する。本発明を適用した装置は、例えばビデオテープレコーダ（VTR）等のハードウェアにより構成してもよく、あるいは、例えば電子番組ガイド（electronic programming guide）等のソフトウェアでもよい。

【0002】

【従来の技術】 欧州特許公開公報EP-A-97118470号には、ホームネットワーク環境において、自らの機能を記述する語彙及び音声インターフェイスを音声装置に送信する機器が開示されている。音声装置は、受信され認識されたユーザの発話に対応するユーザネットワークコマンドに変換し、このユーザネットワークコマンドに基づいて、この機器を制御する。

【0003】 1998年9月に、モトローラ社（Motorola）は、拡張可能マークアップ言語XML（Extensible Markup Language：以下、XMLという。）に基づくV

連する言語を公開した。この言語は、プロンプト及び一般的に利用可能なオプションのリストからなる対話ステップを指定することによって、ダイアログ（対話）システムを記述するのに用いられる。ハイパーメディア記述言語HTML（Hypertext Markup Language：以下、HTMLという。）により、リッチテキストとともに、画像、ハイパーテキストリンク、グラフィックユーザインターフェイス入力コントロールを容易に記述できるように、VoxML 1.0によれば、音声アプリケーションの記述が容易となる。このVoxML 1.0は、複数の要素とこれら要素の属性のシンタクス、用例、Voxマークアップ言語文書又はダイアログの構成、及びVoxマークアップ言語を用いるアプリケーションを開発するときに役立つ他のリファレンスを指示するポインターに関する情報を含む。

【0004】同様に、ギリシャのロートス島における、欧州音響学会主催のユーロスピーチ97で、リチャード・スプロートらによって発表された論文「テキスト音声合成のためのマークアップ言語、ISSN 1018-4074の1774ページ（“a markup language for text-to-speech synthesis” by Richard Sproat et al. ESCA. Eurospeech 97, Rhodes, Greece, ISSN 1018-4074, page 1774）には、音声テキストマークアップ言語（spoken text markup language）（STML）が開示されている。このSTMLは、テキスト音声（TTS）合成器（text-to-speech synthesizers）にテキストの構成に関する知識を提供するものである。STMLテキストは、例えば、多言語のTTSシステムにおいて、言語を設定し、その言語に対して話者を初期化して、適切な言語と話者特定テーブルがロードされるようにする。

【0005】さらに、フィリップス社（Phillips）は、HDDLと呼ばれる対話記述言語を開発した。HDDLは特に、自動問い合わせシステムにおいて用いられる対話のための言語である。対話システムは、HDDLを用いてオフラインモードで作成された後、販売される。

【0006】

【発明が解決しようとする課題】図3は、従来の音声装置21を示す図である。音声装置21は、マイクロホン22とスピーカ30及びバス40に接続されている。マイクロホン22の入力信号は、メモリ23aを内蔵したデジタルシグナルプロセッサ（DSP）23で処理された後、中央演算処理装置（CPU）24に供給される。CPU24は、演算結果を、メモリ29aを内蔵したDSP29を介してスピーカ30に出力し、あるいは、リンクレイヤ制御部25及びI/F物理レイヤ部26を介してバス40に出力する。DSP23及びCPU24は、メモリ28にアクセスすることもできる。メモリ28には、入力信号の処理を制御するために必要な全ての情報が格納されている。さらに、DSP23は、メモリ27の特徴抽出部27fにアクセスする。メモリ27に

は、音声認識及び音声合成に必要な全ての情報が格納されている。CPU24は、複数の機器、ここでは機器#1と機器#2及び機器#3のための各音声インターフェイス定義部及び汎用音声インターフェイス定義部だけでなく、認識部及び書記素／音素変換部27fにもアクセスする。なお、各音声インターフェイス定義部及び汎用音声インターフェイス定義部は、メモリ27内に別々に記憶されている。

【0007】このように従来の音声装置21では、音声インターフェイス定義部を機器別に設けるとともに、汎用音声インターフェイス定義部を別個に備えていたため、複数の機器が共有するコマンドを個別に管理する必要があり、構造が複雑となり、処理や分析の効率が悪かった。

【0008】そこで本発明は、上述した実情に鑑みてなされたものであり、本発明の目的は、ネットワーク機器の機能及び音声インターフェイスをネットワークの音声装置に送信し、音声装置内で複数のネットワーク機器の機能及び音声インターフェイスを操作する簡易且つ効率的なネットワーク機器の制御方法を提供することを目的とする。

【0009】すなわち、本発明の目的は、ネットワーク内で、音声装置によってネットワークに接続されたネットワーク機器を制御するための簡単で、迅速で柔軟性のある方法を提供することである。音声装置とは、ユーザコマンドをユーザネットワークコマンドに変換して、ネットワーク機器の機能及び音声インターフェイスに基づいてネットワークを介してネットワーク機器を制御する装置である。

【0010】

【課題を解決するための手段】本発明によると、ネットワークに接続された全てのネットワーク機器は、ネットワーク機器の機能及び音声インターフェイスを定義する少なくとも1つの機器文書に対応している。機器文書は、ネットワーク機器自身が備えていてもよく、あるいは、ネットワーク機器の外部に作成してもよい。機器文書は、ネットワーク機器のユーザコマンド解釈要素と関連するユーザネットワークコマンドとの対を1組又は複数組有する。さらに、音声装置は、機器文書を受信し、受信した機器文書を統合して1つの音声インターフェイス記述を作成する。この統合された音声インターフェイス記述は、音声装置の言語に基づいて作成される。また、音声インターフェイス記述を、全ネットワークのための汎用文書として参照することができる。音声装置によって受信され認識されたユーザの発話コマンドは、ユーザコマンド解釈要素と汎用言語に含まれる関連するユーザネットワークコマンドの全ての対に基づいて、ユーザネットワークコマンドに変換される。ユーザコマンド解釈及び実行要素は、例えば、語彙要素、文法の定義、発音要素を含む。

【0011】欧州特許公報第E P-A-9711847 0号に開示されているように、音声装置は、機器文書を音声装置が備える記憶装置又はリモートに存在するデータベースからフェッチする。

【0012】汎用文書は、実行時に、ネットワーク機器から機器文書が受信された後、音声装置内で、純粋に統語的に適応化され得る。ネットワーク機器は複数の機器文書を有してもよい。それらの文書それぞれは、ネットワーク機器の機能の一部及びインターフェイスの一部を記述したものである。そして、それらのうちで実際に必要とされる文書だけが音声装置に転送される。このような文書は、例えば、ある言語を用いてネットワーク機器の機能を記述した文書又はネットワーク機器の機能の一部を定義する文書である。音声装置又は個々のネットワーク機器自身は、それら個々のネットワーク機器のさらなる音声機能を必要とすると判定した場合は、そのさらなる音声機能に対応する機器文書を音声装置に送信できる。音声装置は、このさらなる文書に基づいて汎用文書を適応化し、実行時に、その適応化された汎用文書に基づいて対応するユーザネットワークコマンドを生成する。

【0013】

【発明の実施の形態】図1は、本発明を適用した音声装置1を示す図である。音声装置1は、マイクロホン2と、スピーカ10、バス20に接続されている。マイクロホン2に入力された音声に基づく入力信号は、メモリ3aを内蔵したデジタルシグナルプロセッサ(DSP)3で処理された後、中央演算処理装置(CPU)4に供給される。CPU4は、演算結果を、メモリ9aを内蔵したDSP9を介してスピーカ10に出力し、あるいは、リンクレイヤ制御部5及びI/F物理レイヤ部6を介してバス20に出力する。DSP3及びCPU4は、メモリ8にアクセスすることもできる。メモリ8には、入力信号の処理を制御するために必要な全ての情報が格納されている。さらに、DSP3は、メモリ7の特徴抽出部7eにアクセスする。メモリ7には、音声認識及び音声合成に必要な全ての情報が格納されている。

【0014】従来の音声装置では、接続された複数のネットワーク機器のそれぞれに対応する複数の音声インターフェイス定義部と、1つの汎用音声インターフェイス定義部とが設けられていた。一方、本発明を適用した音声装置1は、唯一の統合された音声インターフェイス定義部を有する。この統合された音声インターフェイス定義部は、従来の汎用音声インターフェイス定義部に記述された汎用文書に対応している。

【0015】図2は、本発明を適用したネットワーク機器11の機能的ブロック図である。ネットワーク機器11は、CPU12を備える。CPU12は、ネットワーク機器11を制御するための機器制御用ソフトウェア15と、メモリ14と、リンク層制御部17と、I/F物理

層部16とインタラクションし、種々の情報をバス10に出力する。さらに、CPU12は、メモリ13とインタラクションすることもできる。本発明によると、メモリ13は、音声インターフェイス定義すなわち、1又は複数の機器文書を記憶している。

【0016】上述のように、ネットワーク機器11は、音声インターフェイス定義部を有するメモリ13を必ずしも備えていなくてもよい。ネットワーク機器11の音声インターフェイス定義は、音声インターフェイス定義提供装置により提供してもよい。音声インターフェイス定義提供装置には、図1に示す音声装置1がアクセスすることができる。

【0017】以下に説明する本発明の実施例においては、1つの音声装置1を備えるネットワークに2つのネットワーク機器11が接続されている。初期的には、音声装置1の汎用文書は空で、2つのネットワーク機器11の言語を定義する2つの機器文書は、音声装置1内で1つのインターフェイスの記述に併合(マージ)されている。なお、ネットワークに接続されるネットワーク機器11の数は、2以上の任意のn個であってもよい。また、音声装置1の汎用文書を新たに受信された機器文書に基づいて適応化することができる。以下では、説明を簡潔に行うために、機器文書は、1つの語彙要素からなるユーザコマンド解釈要素を有するものとする。

【0018】以下の説明において、L1は、受理言語(accepted language)すなわち、第1の機器の語彙要素及び対応するコマンドを示す。同様にL2は、第2の機器の語彙要素及び対応するコマンドを示す。数学的表現によれば、L1は、少なくとも1つの語彙要素、すなわち単語 w_i とこの単語 w_i に対応するユーザネットワークコマンドの集合である。なお、単語 w_i は、単一の語ではなくてもよく、複数の単語からなる完全な発話であってもよい。受理言語は、語彙要素に加えて、例えば、発音、単語列のための文法及び/又は音声理解(speech understanding)及び対話のための規則を含んでもよい。

【0019】第1の実施例においては、L1とL2は、同じ語彙要素すなわち共通の単語を含んでいない。したがって、 $L1 \cap L2 = \{\}$ であり、インターフェイス記述のためにマージされた受理言語Lは、 $L = L1 \cup L2$ となる。すなわち、音声装置1内の汎用文書は、第1の機器の言語L1により記述された機器文書1と、第2の機器の言語L2により記述された機器文書2から得られる語彙要素と対応するコマンドの対を併合することにより構築される。L1とL2は同じ語彙要素を含まないので、語彙要素は対応するコマンドとともにその語彙要素がどの機器に対してのものであるかを暗示するため、ユーザネットワークコマンドは、適切に生成されて、対応する機器に正しく送信される。

【0020】第1の実施例において、2つのネットワー

ク機器は、例えばテレビジョンセットとCDプレーヤである。このとき、テレビジョンセットに対応するL1とCDプレーヤに対応するL2はそれぞれ、機器文書内で以下の語彙要素からなる。

【0021】L1 = {MTV, CNN}

L2 = {再生, 停止}

L1とL2は、同じ語彙要素を含まない、すなわちL1 ∩ L2 = {} であるので、インターフェイス記述のマジされた受理言語Lは、L = L1 ∪ L2 = {MTV, CNN, 再生, 停止} となる。例えば、これらの語彙要素はそれぞれ以下のリストに提示されるような機能を有するユーザネットワークコマンドと対応している。

【0022】

MTV → <テレビジョンをMTVに切り換える>

CNN → <テレビジョンをCNNに切り換える>

再生 → <CDプレーヤを再生モードにする>

停止 → <CDプレーヤを停止モードにする>

しかしながら、2つの機器が同じ語彙要素すなわち同じ語を共有する場合、L1 ∩ L2 ≠ {} となる。本発明によれば、これら共通の語彙も識別可能である。本発明の第1の実施例において、機器の名称が各受理言語間、すなわちL1とL2間の少なくとも共通部分を形成する同じ語彙要素に前置又は後置される。したがって、ユーザは、ユーザの発話コマンドを所望の各機器の名称の前に及び／又は後ろに付加しなければならない。上述のように、コマンドが共通のものでない場合は、機器の名称を付加する必要はないが、付加してもよい。

【0023】以下の第2の実施例においては、機器の名称は各コマンドの前に付加され、インターフェイスの記述を形成する新しい言語Lは、2つの言語L1とL2の間で混同されるおそれのない単語の併合と機器の名称が前置された機器の言語で示される。機器の名称をそれぞれn1及びn2とすると、

$$L = L1 \setminus (L1 \cap L2) \cup L2 \setminus (L1 \cap L2) \cup n1L1 \cup n2L2$$

となる。

【0024】以下では、CD再生装置及びテープ再生装置を備えるネットワークを例に、上述の実施例を説明し、本発明によるネットワーク機器の制御方法を明らかにする。ここで、CD再生装置及びテープ再生装置はそれぞれ、「CD」、「テープ」と命名されている。CD再生装置の受理言語L1とテープ再生装置の受理言語L2をそれぞれ構成する語彙要素を以下に示す。

【0025】L1 = {再生, 停止, スキップ}

L2 = {再生, 録音, 停止, 巻き戻し}

音声インターフェイスの受理言語は、 $L = (L1 \setminus (L1 \cap L2)) \cup (L2 \setminus (L1 \cap L2)) \cup n1L1 \cup n2L2$ となる。対応するユーザネットワークコマンド

の機能は、以下のようになる。

【0026】

スキップ → <CDの次のトラックへスキップする>

録音 → <テープ再生装置を録音モードにする>

巻き戻し → <テープを巻き戻す>

CD再生 → <CD再生装置を再生モードにする>

CD停止 → <CD再生装置の再生を停止する>

CDスキップ → <CDの次のトラックへスキップする>

テープ再生 → <テープ再生装置を再生モードにする>

テープ録音 → <テープ再生装置を録音モードにする>

テープ停止 → <テープ再生装置の再生／録音を停止する>

テープ巻き戻し → <テープを巻き戻す>

【0027】この第2の実施例では、同一単語問題 (same words problem) は、認識されたコマンドが多義的であるとき、ユーザに対して自動的に問い合わせを行い、ユーザに所望する機器を明示させることによって解決される。この場合のシナリオは、形式的には、第1の実施例における識別性を有する言語と同じであるが、その解釈が変わる。この解釈のバリエーションを明示するこの実施例を明確にする第3の実施例は、第2の実施例に基づいており、すなわち、ネットワーク機器としてのCD再生装置及びテープ再生装置に付与されたシナリオに基づいている。受理言語Lは、受理言語L1と受理言語L2の共有部分に含まれる語彙要素を認識するときを選択処理が行われるという点では、第1の実施例の受理言語Lと同様に形成される。すなわち、 $L = L1 \cup L2$ である。第2の実施例のように、受理言語L1と受理言語L2がそれぞれ同じ語彙要素を有するという条件のもとで、音声装置の受理言語Lを構成する語彙要素は、 $L = \{\text{再生, 停止, スキップ, 録音, 巻き戻し}\}$ となる。この場合の対応するユーザネットワークコマンドは、

再生 → <明確にせよ # 1>

停止 → <明確にせよ # 2>

スキップ → <CDの次のトラックへスキップする>

録音 → <テープ再生装置を録音モードにする>

巻き戻し → <テープを巻き戻す>

である。ここで、「明確にせよ # 1」は、音声認識部がどの機器に対するコマンドか ("on what device") という問い合わせを出力して、受理言語Lを受理言語Laに切り換えるということの意味する。Laは、 $La = \{\text{CD, テープ}\}$ という語彙要素からなり、対応するユーザネットワークコマンドの機能は、CD → <CD再生装置を再生モードにする> 及び L にスイッチバックする> となり、テープ → <テープ再生装置を再生モードにする> 及び L にスイッチバックする> となる。

【0028】またここで、「明確にせよ # 2」は、音声装置1がどの機器に対するコマンドかという問い合わせを出力して、受理言語Lを受理言語Lbに切り換えるということの意味する。Lbは、Lb = {テープ, CD} となる。

ープ}と示される語彙要素からなり、対応するユーザネットワークコマンドの機能は、CD→<CD再生装置を停止モードにする>及び<Lにスイッチバックする>となり、Tape→<テープ再生装置を停止モードにする>及び<Lにスイッチバックする>となる。

【0029】本発明の第3の実施例において、第2の実施例に示される音声装置の制御処理を短縮することもできる。この場合、音声装置1は直前に用いられた機器をデフォルト機器とみなして、音声装置1とデフォルト機器の統合された受理言語Lからなる汎用文書に基づいて、受信され認識されたユーザの発話コマンドに対応するユーザネットワークコマンドに変換する。直前に用いられた機器以外の他の機器を再度制御する場合にも、音声装置の制御処理を短縮することができる。この場合、例えば単に機器の名称を発することによって、機器を制御するよう音声装置に指示する。

【0030】第4の実施例では、ダイアログの適応化が行われる。上述のような状況においては、多義的な語を正しく割り当てるために、機器明示の問い合わせが、ユーザに対して出力され、システムはその問い合わせに対するユーザの反応を認識する。機器明示のための問い合わせに対するユーザの回答の大部分が直前に指定された機器である場合、問い合わせをスキップし、ユーザネットワークコマンドを直前の機器に直接供給するようにしてもよい。また、コマンドがどの機器に対して出されたのかを予測するために、複数の情報源を用いてもよい。例えば、直前に使用された機器を考慮するだけでなく、各機器の類似点又は機器の分類も考慮する。例えば、記録コマンドは、ユーザがテレビを観ている間に出される場合は、オーディオテープへの記録ではなく、ビデオテープへの記録である可能性が高い。一方、ユーザがラジオを聴いている間に出される記録コマンドは、オーディオテープへの記録である可能性が高い。そこで、使用される可能性の高いネットワーク機器を、例えば電源が投入されているネットワーク機器から選択するようにしてもよい。また、機器の分類は、ユーザの行為又は機器の機能に関する記述を調査することによって学習することもできる。

【0031】上述の全ての実施例において、ネットワーク機器がネットワークに接続された後又は、機器がユーザネットワークコマンドを音声装置1から受信した後に、各ネットワーク機器は、機器文書を音声装置1に直接送信することができる。また、現在のステータスから独立した機器文書等を音声装置に送信してもよい。機器は、複数の機器文書を有していてもよい。または、機器は、内容情報(content information)の変更に基づいて動的に機能を変えて、その変化を認めるときにダイアログ及び音声認識能力を更新してもよい。あるいは、新しい文書を音声装置に送信することによって、定期的に機能を変えて、機器文書は、対応する機器内に格納

されずに、例えば装置製造者のインターネットサーバのような機器文書提供装置内に格納されてもよい。この場合、音声装置1は上記の機器がネットワークに接続されていると認識すると、機器文書をダウンロードしてもよい。

【0032】機器は、それぞれが異なる国の言語による機器全体の受理言語を含む複数の文書を有することができる。例えば、このような機器は、最初にドイツ語の機器文書を音声装置に送信する。続いて、ユーザから対応するコマンドを受信した後に、英語の機器文書を音声装置に送信することができる。これにより、音声装置は、ドイツ語又は英語のいずれかのユーザコマンドをユーザネットワークコマンドに変換して機器を制御することができる。

【0033】語彙要素の他に、認識又は合成を行うときに用いられる語のための発音の定義を機器文書の要素として含ませてもよい。そして、それらの発音の定義を音声装置の汎用文書に直接統合してもよい。発音の定義が言語に含まれるとき、それらの発音の定義は、例えば、1つの機器用の認識及び合成部に専用にもよく、又は、全ての機器で共有してもよい。すなわち、第1の機器によって導入された発音の定義は、第2の機器、すなわちネットワークに接続された他の機器に対して発せられた語に対しても有効とすることができる。

【0034】合成を行うために、複数の機器文書を統合し、音声装置内で、汎用文書の適応化を行ってもよい。問い合わせに対する回答を1つの機器だけでなく、複数の機器に対して行ってもよい。複数の機器からの発話の合成は、同時に出力されないことを保証する必要がある。発話の中には、供給優先度に基づいて、例えば警告のような外部の事象によって生じるものがある。例えばそれらの発話は、他の発話に割り込み、ユーザの非活動期間が終わると、割り込まれた発話は再導入され、ダイアログは続行する。ユーザの入力及びユーザの発話に一致する尤度が高いと考えられるストリームが継続している間、全てのストリームがアクティブな状態にあるという事実を前提として、複数のダイアログのストリームが平行して管理される。

【0035】音声装置の汎用文書は、初期的には空であってもよく、あるいは、1又は複数のネットワーク機器のための、インターフェイスを記述する要素の基本セットを有していてもよい。後者の場合、機器から音声装置に送信される文書は、具体性を有さず、汎用文書を部分的に参照するものであってもよい。

【0036】例えば、特定の機器のためのダイアログ文法のみを機器文書において特定し、発話の言い回しのためのデフォルト文法を汎用文書のユーザコマンド解釈要素の基本的なセットに含ませてもよい。また、音声認識部の汎用文書内に、ある機能を完全に記述してもよい。例えば、音声認識ガイドを、音声ダイアログにトピック

どのように制御するのかを記述する。この場合、デフォルト文法及び一連の単語のための文法には、例えば、番組名、製作者名、俳優名、時間帯及び放送日に関する情報を示す変数が含まれており、電子番組ガイド装置から音声装置に送信される文書は、これら変数を満たす情報のみを含む。

【0037】他の実施例においては、機器から音声装置に送信される機器文書は、ネットワークコマンドに対応するキーワードとカテゴリ識別子を含む。音声装置は、これらのカテゴリ識別子に基づいて、発生し得る連続するユーザの発話に対して、単語列に関するどの文法を適用するかを決定する。例えば、キーワードが「ボイス・オブ・アメリカ (voice of America)」であり、カテゴリが「ラジオ」である場合、ユーザは、「ボイス・オブ・アメリカを聴きたい (I want to listen to voice of America)」又は「ボイス・オブ・アメリカをオン (Turn on voice of America)」等と発生する。また、カテゴリが「ラジオカセットレコーダ」である場合、発生し得るユーザの発話は、「ボイス・オブ・アメリカを録音 (Please record voice of America)」等である。

【0038】汎用文書の適応化を行うために、音声理解の基礎を形成する概念／値の組、例えばある処理を開始するために必要な(時刻, 8:00)といった組を音声装置に送信される機器文書に機器用の文法として含ませてもよい。音声インターフェイスを記述する機器文書は、規則を含み、この規則は、単語列と概念との間のマッピング及びどのように変数値を埋めるかを定義する。さらに、機器文書は、概念／値を一对の組から動作へのマッピングを定義する規則を有する。例えば(コマンド, 記録)、(タイトル, 風とともに去りぬ)及び(時刻, 8:00)という組は、VTRの録画処理を表す。ダイアログは、ユーザ側からはじめてもよく、すなわち発話の順番は、ユーザを主導としてもよい。発話の順番とは、例えば、「8時、映画、風とともに去りぬ」等である。1又は複数の発話で情報が伝えられる場合は、発話の順番は、ユーザに任せられる。システムは、例えば、文法に基づきそして各ネットワーク機器に送信されたユーザネットワークコマンドによって動作を開始するために必要な情報を収集する。このシステムは、例えば、付加的な文法を含むさらなる文書を音声装置に送信し、音声装置は、情報が欠落している場合、ユーザに情報を聞き返すことができ、また、新しく認識されたユーザの発話をさらなるユーザネットワークコマンドに割り当てることができる。

【0039】機器文書は、ユーザコマンド解釈要素として、語彙要素、発音、単語列のための文法、音声理解及びダイアログのための規則に関する情報の他に、同一又は同等のユーザネットワークコマンドに割り当てられる複数言語に対応する同様な情報、すなわち語彙要素、発音、単語列のための文法、音声理解及びダイアログのた

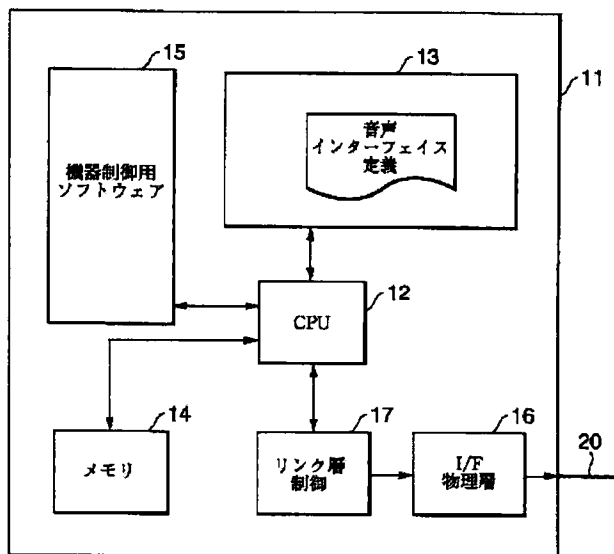
めの規則に関する情報を含んでもよい。この情報によって、ユーザはいかなる言語によっても機器を制御することができる。上述のように、この情報を同じ機器が有する複数の機器文書中に含ませてもよい。これによって、複数言語の音声による処理を実現することができる1つの音声装置のインターフェイスを特定できる。上記の概念／値の対は言語独立であるので、ダイアログの文法は言語独立とすることができる。

【0040】上述のように、機器文書は、発話のための文法によって明示的にあるいは暗示的に付与された語彙を含む。1又は複数の発音を各語彙要素、すなわち各語に対して任意に付与することもできる。任意の発音が機器文書に付与されないときは、これを自動的に生成してもよい。しかしこの場合、特に単語が固有名詞、外来語、略語である場合や方言又は外国語による特殊な発音を認識しなくてはならない場合、エラーが発生しやすい。発音は、単音又は異音のモデルのセットが与えられた認識部内に単語のモデルを生成するための基礎となる。

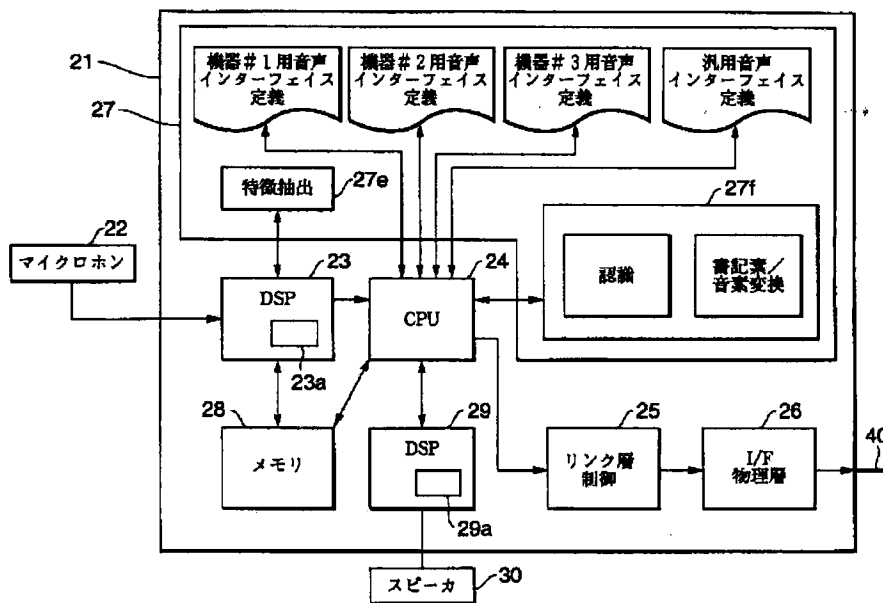
【0041】機器文書は、複数の利用可能な機器に関する複合コマンドを作成するユーザコマンド解釈要素及び関連するユーザコマンドの他に、標準化された情報、例えば「VTR」のような機器のカテゴリに関する情報及び／又は「記録可能なテープを有する」等の機能に関する情報を含んでもよい。例えばユーザは、テレビジョンセットに表示される電子番組ガイドから将来放映される予定の映画を選択するとき、「この映画を録画して下さい (Please record this movie)」と発する。つまり、ネットワークコマンドが、適切なチャンネル、日付及び時間帯に関する情報を番組ガイド装置からVTRに送信する必要がある。そして、適切な録画が行われることを保証するコマンドをVTRに送信する必要がある。これに類似した状況として、映画のあるVTRテープから別のVTRテープにダビングする場合がある。このような状況では、各機器の音声インターフェイス記述は、通常、このような発話、すなわち適切なチャンネル、日付及び時間帯等に関する発話を定義しない。したがって、音声装置は、推論部による推論に基づいて汎用文書を適応化してもよい。推論部は、推論構成要素として機器文書内に含まれており、単語列に関する文法及びダイアログを含む音声インターフェイスの記述を複数の機器に関する複合コマンドの1つに含まれる可能性のある追加的機能に提供する。

【0042】ネットワーク機器を制御するための全ての情報は、音声装置内の1つの汎用文書に含まれるので、処理あるいは分析を容易に行うことができる。特に、ユーザコマンド解釈要素の異なる利用可能な内容間又は／及び汎用文書内の同様又は同等のユーザネットワークコマンドに割り当てられた二重の情報間に関連があるときは特に空目である。

【図2】



【図3】



フロントページの続き

(51) Int. Cl. 7

H04Q 9/00

識別記号

331

F I

G10L 3/00

テームコード (参考)

561H

(72)発明者 ステファン ラップ
ドイツ連邦共和国 ディー－70736 フェ
ルバッハシュトゥットゥガルトー シュト
ラーセ 106 ソニー インターナシヨナ
ル (ヨーロッパ) ゲゼルシャフト ミッ
ト ベシュレンクテル ハフツング シュ
トゥットゥガルト テクノロジーセンター
内

(72)発明者 シルケ ゴロンジー
ドイツ連邦共和国 ディー－70736 フェ
ルバッハシュトゥットゥガルトー シュト
ラーセ 106 ソニー インターナシヨナ
ル (ヨーロッパ) ゲゼルシャフト ミッ
ト ベシュレンクテル ハフツング シュ
トゥットゥガルト テクノロジーセンター
内

(72)発明者 ラルフ コンペ
ドイツ連邦共和国 ディー－70736 フェ
ルバッハシュトゥットゥガルトー シュト
ラーセ 106 ソニー インターナシヨナ
ル (ヨーロッパ) ゲゼルシャフト ミッ
ト ベシュレンクテル ハフツング シュ
トゥットゥガルト テクノロジーセンター
内

(72)発明者 ペーター ブフナー
ドイツ連邦共和国 ディー－70736 フェ
ルバッハシュトゥットゥガルトー シュト
ラーセ 106 ソニー インターナシヨナ
ル (ヨーロッパ) ゲゼルシャフト ミッ
ト ベシュレンクテル ハフツング シュ
トゥットゥガルト テクノロジーセンター
内

(72)発明者 フランク ジロン
ドイツ連邦共和国 ディー－70736 フェ
ルバッハシュトゥットゥガルトー シュト
ラーセ 106 ソニー インターナシヨナ
ル (ヨーロッパ) ゲゼルシャフト ミッ
ト ベシュレンクテル ハフツング シュ
トゥットゥガルト テクノロジーセンター
内

(72)発明者 ヘルムート ルッケ
ドイツ連邦共和国 ディー－70736 フェ
ルバッハシュトゥットゥガルトー シュト
ラーセ 106 ソニー インターナシヨナ
ル (ヨーロッパ) ゲゼルシャフト ミッ
ト ベシュレンクテル ハフツング シュ
トゥットゥガルト テクノロジーセンター
内